

Edit distance between unlabeled ordered trees

Anne Micheli

Dominique Rossin

October 31, 2018

Abstract

There exists a bijection between one stack sortable permutations –permutations which avoid the pattern 231– and planar trees. We define an edit distance between permutations which is coherent with the standard edit distance between trees. This one-to-one correspondence yields a polynomial algorithm for the subpermutation problem for (231) avoiding permutations.

Moreover, we obtain the generating function of the edit distance between ordered trees and some special ones. For the general case we show that the mean edit distance between a planar tree and all other planar trees is at least $n/\ln(n)$.

Some results can be extended to labeled trees considering colored Dyck paths or equivalently colored one stack sortable permutations.

1 Introduction

The edit distance between two trees is the minimal number of edit operations to transform one tree into the other. The edit operations are deletion (edge contraction), insertion of an edge and relabeling of a vertex.

The main problem is to find efficient algorithms to compute this distance between ordered labeled trees. Many algorithms have been proposed [1, 2]. The basic idea of all these dynamic algorithms arises from the paper of Zhang and Shasha [1]. Further improvements have been made [2].

Comparing the structure of molecules and finding the preserved ones during a genetic mutation can be seen as an edit distance problem. The application field of this problem is not restricted to biology: in computer vision, objects are represented by their skeletons -which are trees-, and in computer science, edit distance is used to compare structural similarities between XML documents [3].

But no combinatorial interpretation has been made of the edit distance between trees. In this article, we introduce one-stack sortable permutations [4, 5]. These one-stack sortable permutations are (231) pattern-avoiding permutations and we show that they are in one-to-one correspondence with ordered trees.

Moreover the edit operations can be easily described in terms of one-stack sortable permutations. This leads to a purely combinatorial explanation of the edit distance.

Some polynomial algorithms are known to compute the edit distance between trees [1]. By our correspondence, we show that computing the greatest common pattern between two (231)-avoiding permutations is also polynomial whereas it is NP-complete for general permutations [6].

2 Definitions

2.1 One-stack sortable permutations

We describe in this section an encoding for planar trees. We number the edges of the tree by a postfix traversal and then read the permutation by a prefix traversal. The obtained permutations are called one stack sortable permutations [4, 5]. An alternate definition is the following:

Definition 1. Let $n \in \mathbb{N}$, a one-stack sortable permutation on $\{1 \dots n\}$ is a permutation σ such that $\sigma = InJ$ where I and J are one-stack sortable permutations on $\{1 \dots p\}$ and $\{p+1 \dots n-1\}$ respectively. Notice that I or J could be empty.

Note that in the sequel, permutations are seen as words.

Theorem 1. One-stack sortable permutations are in one-to-one correspondence with rooted ordered trees.

Proof. Given a tree T with n edges, number the edges by a postfix Depth First Search Traversal (DFS). Read it again by a prefix DFS. It is clear that the obtained permutation is of the form InJ . Moreover I corresponds to the encoding by a postfix DFS of the left subtree as shown in Figure 1. The same goes for J but its numbers are shifted.

Conversely, take a one-stack sortable permutation $\sigma = InJ$.

- If $\sigma = k$ then the corresponding tree is a single edge.
- If $\sigma = InJ$ then the corresponding tree T_σ is the tree obtained by taking an edge $e = (xy)$ (corresponding to n) where x is the root of T_σ . Since I and J are also one-stack sortable permutations, we can recursively build the corresponding trees T_I and T_J . Put them at each end of the edge e , ie T_I is hanging on x such e is the rightmost edge of x , and T_J on y .

This construction is unique.

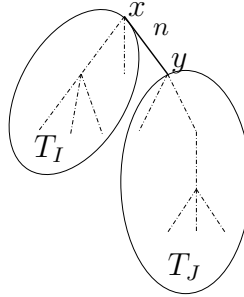


Figure 1: Coding a tree with a one-stack sortable permutation.

□

If σ is a one-stack sortable permutation, let $\mathcal{T}(\sigma)$ denote the tree associated to σ . Conversely, if T is a tree, its associated one-stack sortable permutation is denoted by $\Theta(T)$. Moreover, in the sequel, σ_k will either denote the k -th letter of the word σ or the corresponding edge in $\mathcal{T}(\sigma)$.

Definition 2. A subsequence of a permutation $\sigma = \sigma_1 \dots \sigma_n$ is a word $\sigma' = \sigma_{i_1} \dots \sigma_{i_k}$ where i_1, \dots, i_k is an increasing sequence of elements of $\{1, \dots, n\}$.

Let Φ be the bijective mapping of $\{\sigma_{i_1}, \sigma_{i_2}, \dots, \sigma_{i_k}\}$ on $\{1, \dots, k\}$ preserving the order on σ_{i_l} .

The normalized subsequence (pattern) $\hat{\sigma}'$ is equal to $\Phi(\sigma')$.

Remark 1. The one-stack sortable permutations are the permutations avoiding the normalized subsequence (pattern) 231 [7].

2.2 Edit distance

We briefly recall the definition of the edit distance between trees. Given two trees, the edit distance is the minimal number of operations necessary to transform one into the other. The operations are:

- **Deletion** : This is the contraction of an edge; two vertices are merged. Only one label is kept.

- **Insertion** : This is the converse operation of deletion.

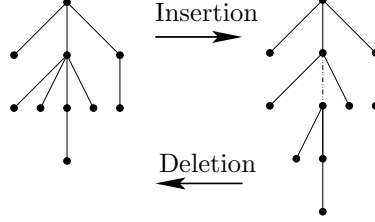


Figure 2: Insertion and Deletion operations on a tree.

A cost can be given to each operation. In this article we take 1 for every cost.

3 Distance on one-stack sortable permutations

Since one-stack sortable permutations are in one-to-one correspondence with planar trees, we define similar edit operations between one-stack sortable permutations and show that these definitions match with edit distance between trees. Moreover, we give a combinatorial interpretation of the distance.

A *factor* of a permutation $\sigma = \sigma_1\sigma_2\ldots\sigma_n$ is a *factor* of the word $\sigma_1\sigma_2\ldots\sigma_n$ i.e.a word of the form $\sigma_k\sigma_{k+1}\ldots\sigma_{k+l}$.

A factor f is *compact* if it is a permutation of an interval of \mathbb{N} .

A factor f of σ is *complete* if no non-empty factor g of σ verifies both:

1. fg is compact where fg is the concatenation of the words f and g ;
2. the greatest element of fg is equal to the greatest element of f .

Take for example the one-stack sortable permutation $\sigma = (1524376)$. The complete factors of σ are $\{1\}, \{15243\}, \{1524376\}, \{5243\}, \{524376\}, \{2\}, \{243\}, \{43\}, \{3\}, \{76\}, \{6\}$.

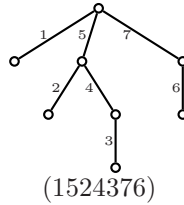


Figure 3: Tree associated to $\sigma = (1524376)$.

A subtree T' of T is a tree such that $T \setminus T'$ is connected.

Lemma 1. *Each compact factor of σ are in one-to-one correspondance with:*

- to a subtree
- to a internal path P in $T = \mathcal{T}(\sigma)$ where each internal vertex of P is of degree 2 in T and P does not end at a leaf (P can be an internal edge).

Proof. First let prove that the subset of edges correponding to a compact factor is connected.

Let σ' be a compact factor of $\sigma = \Theta(T)$. Let $E_{\sigma'}$ be the set of edges corresponding to σ' in T . Suppose that $E_{\sigma'}$ is not connected. Let E_1 and E_2 be two connected components. Let v be the first common ancestor of E_1 and E_2 . Let P_1 (resp. P_2) be the path starting from v and ending at the first vertex of E_1 (resp. E_2). Note that we can choose E_1 and E_2 such that edges of P_1 and P_2 are not in $E_{\sigma'}$. Suppose that P_1 is at the left of P_2 (See Figure 3). In the prefix DFS of T , edges of

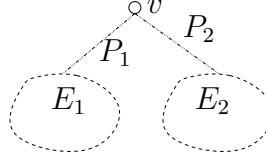


Figure 4: Compact factors are connected components.

P_2 are visited between those of E_1 and E_2 . Thus they should appear in σ' , hence $P_2 = \emptyset$. Thus $v \in E_2$ so that P_1 links E_2 and E_1 . In the postfix DFS, the edges of P_1 have labels greater than those of E_1 and less than E_2 . If $P_2 \neq \emptyset$, it implies that σ' is not compact. Thus $E_{\sigma'}$ is connected.

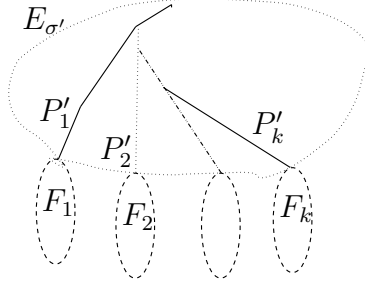


Figure 5: Subtree of T induced by $E_{\sigma'}$.

Consider the subtree T' of T induced by $E_{\sigma'}$. It consists of $E_{\sigma'}$ plus all vertices of T that have an ancestor in $E_{\sigma'}$ as shown in Figure 5.

$E_{\sigma'}$ can be decomposed into edge-disjoint paths P_i thanks to the prefix DFS (See Figure 5). F_i is the subtree pending on P_i which can be empty.

The prefix DFS of T' (which is a factor of σ) gives the associated permutation $\Theta(P'_1)\Theta(F_1)\Theta(P'_2)\Theta(F_2)\dots\Theta(P'_k)\Theta(F_k)$. So $\sigma' = \Theta(P'_1)\Theta(F_1)\Theta(P'_2)\Theta(F_2)\dots\Theta(P'_k)$, hence $F_i = \emptyset, \forall i < k$.

- Suppose $F_k \neq \emptyset$. If $k > 1$, then the edges of F_k are visited after at least one edge of P'_1 , and before the edges of P'_k in the postfix DFS. Since σ' is compact, it implies $k = 1$.
- If $F_k = \emptyset$, $E_{\sigma'}$ is a subtree.

The converse is straightforward.

□

Proposition 1. *The set of complete factors of σ corresponds to the set of subtrees of the associated tree.*

Proof. Let T' be a subtree of T and $\sigma = \Theta(T)$. The edges of T' are visited consecutively by the postfix (resp. prefix) DFS of T . Thus the sequence of edges of T' is a compact factor $\sigma_k\sigma_{k+1}\dots\sigma_{k+l}$ of σ . σ_{k+l+1} is an edge which is visited after all edges of T' by the prefix DFS. Thus it is the first time this edge is visited by the traversal. Hence, its label is greater than those of T' . Thus $\sigma_k\sigma_{k+1}\dots\sigma_{k+l}$ is complete.

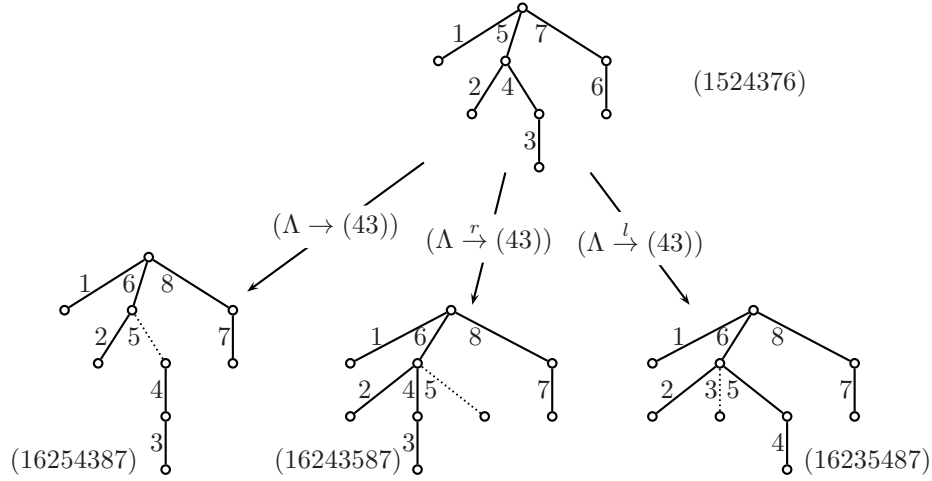


Figure 6: Insertion operations for $f = (43)$.

Conversely, let σ' be a complete factor. As σ' is compact, by Lemma 1, it corresponds either to a subtree or to an internal path P with a subtree F hanging on P . $\Theta(P)\Theta(F) = \sigma'\Theta(F)$ is also a compact factor of σ and it has the same maximum as σ' which contradicts the completeness of σ' . \square

Remark 2. Let σ be a one-stack sortable permutation and $\sigma_k = (p(v_k)v_k)$ an edge where $p(v_k)$ denote the parent of v_k . Let σ' be the shortest complete factor of σ such that $\sigma' = \sigma_k\sigma_{k+1}\dots\sigma_{k+l}$ where $\sigma_i = (p(v_i)v_i)$. By previous proposition $\mathcal{T}(\sigma')$ is a subtree of $\mathcal{T}(\sigma)$. The children of v_k are the vertices v_{k+i} such that $i \leq l$ and $\sigma_k > \sigma_{k+i} > \sigma_{k+1}, \sigma_{k+2}, \dots, \sigma_{k+i-1}$.

Let $\sigma = \sigma_1 \dots \sigma_k$ be a word of $\{1 \dots n\}$ and a be a letter of $\{1 \dots n\}$. We denote by $[\sigma]_a$ the word $\sigma'_1 \dots \sigma'_k$ where

$$\sigma'_i = \begin{cases} \sigma_i & \text{if } \sigma_i < a \\ \sigma_i + 1 & \text{otherwise} \end{cases}$$

Definition 3. We define two operations on permutations which map the standard definition on trees $([1])$:

1. *Deletion* : Let $1 \leq k \leq n$. The deletion $(\sigma_k \rightarrow \Lambda)$ is the removal of σ_k in a permutation σ and the renormalization on S_{n-1} of the result. We will either talk about the deletion of the edge σ_k or the deletion of the vertex v such that σ_k is the edge $p(v)v$.
2. *Insertion* (see Figure 6) : $(\Lambda \rightarrow \emptyset)$ corresponds to the transformation of the permutation $\sigma = \emptyset$ into $\sigma' = (1)$. If $\sigma \neq \emptyset$, let f be a complete factor of σ . Then, $\sigma = uv$ with u, v factors of σ .
 - (a) $(\Lambda \rightarrow f)$: The resulting permutation is $\sigma' = [u]_a a [v]_a$, $a = \max\{f\} + 1$. This corresponds to the insertion of an inner vertex with $\mathcal{T}(f)$ as subtree.
 - (b) $(\Lambda \xrightarrow{r} f)$: The resulting permutation is $\sigma' = [u]_a f a [v]_a$, $a = \max\{f\} + 1$. This corresponds to the insertion of a leaf as the right sibling of $\mathcal{T}(f)$.
 - (c) $(\Lambda \xrightarrow{l} f)$: The resulting permutation is $\sigma' = [u]_a a [f]_a [v]_a$, $a = \min\{f\}$. This corresponds to the insertion of a leaf as the left sibling of $\mathcal{T}(f)$.

We study now these operations on the permutation $\sigma = (1524376)$.

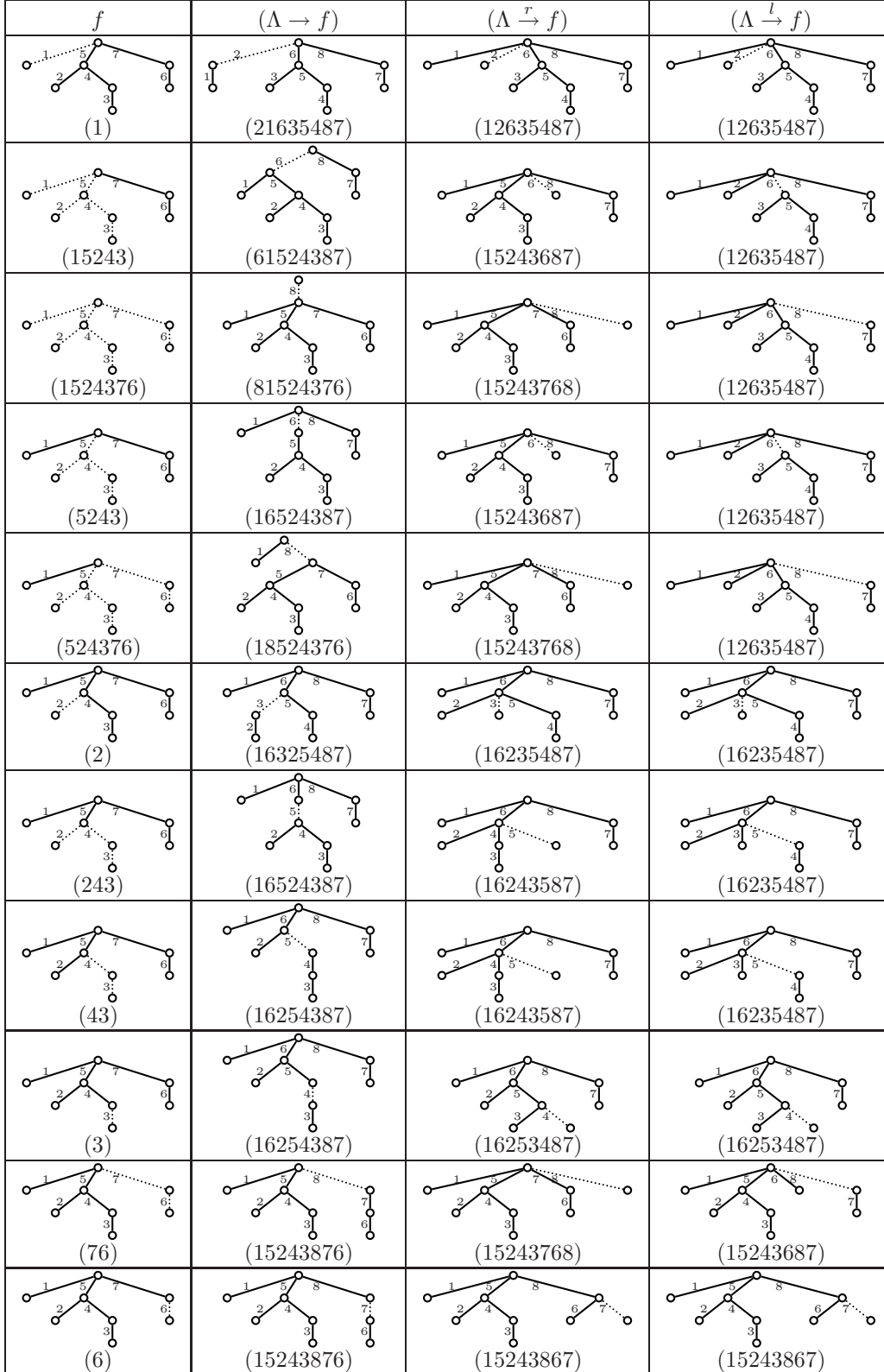


Figure 7: Insertion in permutation $\sigma = 1524376$.

The array of Figure 7 gives all the permutations that can be obtained with a single insertion in σ .

We prove now that the operations (deletion and insertion) defined on one-stack sortable permutations are in fact internal operators for one-stack sortable permutations. Moreover, these operators define an edit distance between permutations coherent with the usual edit distance between trees.

Lemma 2. *The Deletion/Insertion algorithm yields a one-stack sortable permutation.*

Proof. • Deletion : The proof is straightforward considering the one-to-one correspondence with trees and one-stack sortable permutations. Consider a tree labeled by a depth first traversal. Deleting the edge i from this tree changes all labels greater than i by subtracting 1.

• Insertion : Let σ be a one-stack sortable permutation and f be a complete factor of $\sigma = uf$. By Proposition 1, f corresponds to a subtree of $\mathcal{T}(\sigma)$.

1. $(\Lambda \rightarrow f)$: Let $T = \mathcal{T}(\sigma)$ and (e_1, e_2, \dots, e_n) be the edges of T ordered by a prefix DFS of the tree. Note that $\sigma = \alpha_T(e_1)\alpha_T(e_2)\dots\alpha_T(e_n)$ where $\alpha(i)$ is the label of the edge i in T .

Let T' be the tree obtained by the insertion of an internal vertex v ($a = (p(v)v)$) at the root vertex of the subtree $\mathcal{T}(f)$. Moreover $\mathcal{T}(f)$ is a subtree hanging on v . Let $\sigma'' = \Theta(T')$. A prefix traversal of T' orders the edges of T' as follows: $(e_1, e_2, \dots, e_l, a, e_{l+1}, \dots, e_n)$.

Since σ'' is obtained by a prefix traversal, $\sigma'' = u'af'v'$. Since the edges of f appear before a in the postfix DFS, $f' = f$. The edge a in a postfix DFS appears just after f . Thus its label is $\max\{f\} + 1$. All the edges visited after f in T (and so after a in T') by the postfix DFS have their labels increased by 1. Thus $\sigma'' = [u]_a af[v]_a = \sigma'$.

2. $(\Lambda \xrightarrow{l} f), (\Lambda \xrightarrow{r} f)$: The same arguments as for $(\Lambda \rightarrow f)$ hold.

□

Proposition 2. *Insertion and deletion are inverse operations.*

Proof. There are two different kinds of deletions in a tree T .

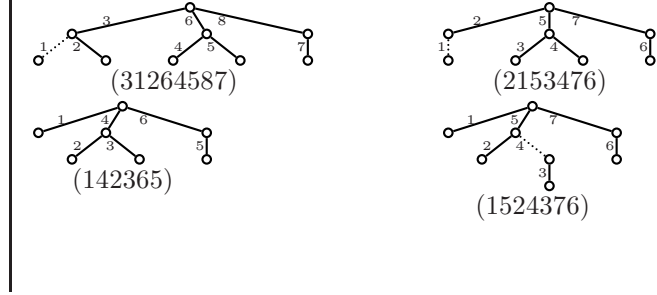
1. Deletion of an inner vertex v . Consider the subtree T' of T hanging on v . It corresponds to a complete factor f in $\sigma = \Theta(T)$. This contraction corresponds to the inverse operation of $(\Lambda \rightarrow f)$.
2. Deletion of a leaf. There are three different cases:
 - Deletion of a vertex with no sibling. This is the same as deleting the parent of this vertex which is an inner vertex except if the tree is reduced to a single edge.
 - Otherwise, this vertex has either:
 - A left sibling v' . Consider the subtree hanging at v' (including $p(v')v'$). It corresponds to the factor f . The inverse operation is $(\Lambda \xrightarrow{r} f)$
 - A right sibling v' . Consider the subtree hanging at v' (including $p(v')v'$). It corresponds to the factor f . The inverse operation is $(\Lambda \xrightarrow{l} f)$

□

Definition 4. *The distance between two one-stack sortable permutations σ_1 and σ_2 is the minimal number of operations -deletion or insertion - to transform σ_1 into σ_2 .*

For example let $\sigma_1 = 31264587$ and $\sigma_2 = 1524376$. We want to transform σ_1 into σ_2 .

- $31264587 \xrightarrow{(1 \rightarrow \Lambda)} 2153476$
- $2153476 \xrightarrow{(1 \rightarrow \Lambda)} 142365$
- $142365 \xrightarrow{(\Lambda \rightarrow 3)} 1524376$



Theorem 2. *The edit distance between ordered trees is the distance between the associated one-stack sortable permutations.*

Proof. This is a consequence of Proposition 2. □

Theorem 3. *The edit distance between one-stack sortable permutations σ_1 and σ_2 is equal to*

$$|\sigma_1| + |\sigma_2| - 2|u|$$

where u is a largest normalized subsequence (pattern) of σ_1 and σ_2 .

Proof. The edit distance $d(\sigma_1, \sigma_2)$ between σ_1 and σ_2 is given by the minimal number of insertions and deletions. If t_1 is an insertion and t_2 is a deletion then there exist a deletion t'_1 and an insertion t'_2 such that $t_1 t_2(\sigma) = t'_1 t'_2(\sigma)$. Note that t'_1 and t'_2 depend on the one-stack sortable permutation σ .

Considering the sequence of edit operations, there exists a sequence made of deletions then insertions that transforms σ_1 into σ_2 . We denote this sequence by $D_1 \dots D_l O_1 \dots O_k$, $l + k = d(\sigma_1, \sigma_2)$.

Consider the one-stack sortable permutation $\sigma' = D_1 \dots D_l(\sigma_1)$. Take $u = \sigma'$. u is a normalized subsequence of σ_1 because deleting an edge from a one-stack sortable permutation yields a normalized subsequence of the original one-stack sortable permutation. u is also a normalized subsequence of σ_2 because inserting an edge in a one-stack sortable permutation s yields a one-stack sortable permutation s' and s is a normalized subsequence of s' .

Conversely, take u as a maximal normalized subsequence of σ_1 and σ_2 . It is straightforward to find $|\sigma_1| - |u|$ operations of deletions such that those deletions transform σ_1 into u . The same goes for σ_2 and u . □

Corollary 1. *Finding the greatest common pattern between two one-stack sortable permutations is polynomial.*

In [6], they proved that finding the greatest common pattern between two permutations is NP-complete. We prove here that the problem becomes polynomial when restricting to one-stack sortable permutations, ie (132) or (231)-avoiding permutations. In fact, the algorithm of Zhang and Shasha [1] on trees solves the problem on one-stack sortable permutations because the algorithm outputs not only the distance but also the greatest common subtree.

4 Lower bounds on average edit distance

In this section we study the average edit distance between a given planar tree T with n vertices and all other planar trees with n vertices. We show that this average distance is lower bounded by $\frac{n}{\ln(n)}$.

Lemma 3. *Let T be a planar tree with n vertices. There are at most $n - 1$ different deletions and $3n^3$ insertions allowed in T .*

Proof. The number of deletions is upper bounded by the number of edges i.e. $n - 1$.

The number of insertions is bounded by 3 times the number of subtrees (or complete factor of the corresponding permutation). The number of subtrees of T rooted at vertex v is bounded by $d(v)^2$ where $d(v)$ denotes the degree of vertex v . Thus the total number of subtrees is bounded by $\sum_v d(v)^2$. \square

Theorem 4. *Let T_0 be a tree with n vertices. The proportion of planar trees with n vertices at distance at most $\mathcal{O}(n/\ln(n))$ tends to 0.*

The average distance between T_0 and the set of planar trees is lower bounded by $n/\ln(n)$.

Proof. Let T_0 be a planar tree. Let $A_k = \{T \in \mathcal{T}_n, \text{dist}(T_0, T) \leq 2k\}$. Note that $A_0 = \{T_0\}$. A tree $T_k \in A_k$ is obtained from T_0 by $l \leq k$ deletions then l insertions. Thus $|A_k| < (n-1)^k (n^3)^k < n^{4k}$. But the number of planar trees $C_n \equiv \frac{4^n}{n\sqrt{\pi n}}$. So that the proportion of planar trees at distance at most $\mathcal{O}(n/\ln(n))$ tends to 0.

Hence the average distance is lower bounded by $n/\ln(n)$. \square

5 Generating functions

Using the combinatorial interpretation of the distance, we compute the generating functions of the edit distance between planar trees with n edges and some special ones as shown in Figure 8. Moreover, we deduce the average distances from the generating functions.

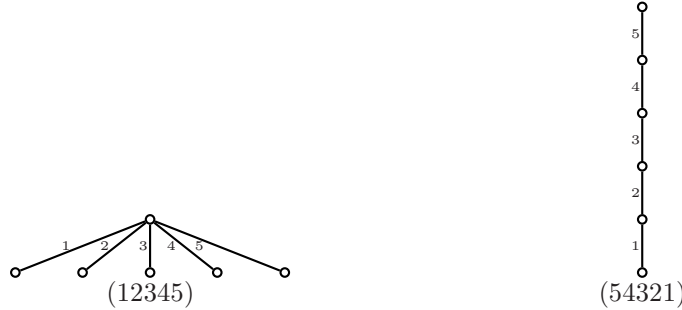


Figure 8: Some canonical trees.

5.1 Generating function of the edit distance between one-stack sortable permutations and $Id = 1\ 2\ \dots\ n$

We denote by $S_1(t, q)$ the generating function of one-stack sortable permutations where t counts the size of the permutation and q the edit distance between one-stack sortable permutations and Id . This is the distance between a tree and the trivial one which is made of n edges and of height 1.

Tree interpretation of the largest increasing subsequence

Proposition 3. *The length of a largest increasing subsequence of a one-stack sortable permutation is the number of leaves of the associated tree.*

Proof. Let T be a planar rooted tree and σ the associated one-stack sortable permutation. We call a leaf-edge an edge incident to a leaf.

1. The subsequence of σ made of the leaf-edges is increasing because the order in which the leaf-edges are visited by a prefix traversal is the same than by a postfix traversal.

2. Suppose that we take an increasing subsequence σ' of σ . This subsequence is in one-to-one correspondence with some edges in the tree. Suppose that there is an internal one $\gamma = (p(\nu)\nu)$. Then, by the postordering of the edges, each edge $(p(v)v)$ such that $\nu = p(v)$ has a smaller label and appears in σ after the edge γ . Thus, none of these edges are in σ' . Moreover, there is at least one leaf edge belonging to the subtree T_γ hanging on ν . Replace edge γ by a leaf of T_γ . The prefix traversal ensures that the obtained subsequence is an increasing one.

□

Proposition 4. *The number of rooted planar trees with n edges and k leaves is equal to the number of rooted planar trees with n edges and $n + 1 - k$ leaves.*

Proof. This is a direct consequence of the symmetry of the Narayana numbers $\frac{1}{n} \binom{n}{k} \binom{n}{k-1}$ which count the number of planar trees with n edges and k leaves.

□

Generating function We now compute the generating function $I(t, p)$ of one-stack sortable permutations of size t and largest increasing subsequence of size p .

- $[I(t, p)]_0 = 1$
- $[I(t, p)]_1 = p$
- $[I(t, p)]_2 = (p + p^2)$

$$[I(t, p)]_n = p[I(t, p)]_{n-1} + \sum_{i=0}^{n-2} [I(t, p)]_i [I(t, p)]_{n-1-i} \quad (1)$$

This formula comes from the decomposition of a one-stack sortable permutation σ into InJ with $n \geq 1$. The largest increasing subsequence of σ is the union of the largest one of I and the largest one of J unless J is empty - in this case, the largest subsequence is the largest one for In .

From this formula we deduce:

$$I(t, p) = 1 + (p - 1)tI(t, p) + tI^2(t, p) \quad (2)$$

- 1 comes from the case $n = 0$ in the equation (1).
- $ptI(t, p)$ comes from $p[I(t, p)]_{n-1}$.

It follows from equation (2):

$$I(t, p) = \frac{1 + (1 - p)t - \sqrt{(p - 1)^2 t^2 - 2(p + 1)t + 1}}{2t} \quad (3)$$

Let $\tilde{S}_1(t, q)$ be the generating function of the difference between the lengths of the one-stack sortable permutation and the largest increasing subsequence in it.

- $[\tilde{S}_1(t, q)]_0 = -1$
- $[\tilde{S}_1(t, q)]_1 = 0$
- $[\tilde{S}_1(t, q)]_2 = q$

Lemma 4.

$$I(t, p) = 1 + p + p\tilde{S}_1(t, p) \quad (4)$$

Proof.

$$\begin{aligned}
I(t, p) &= \sum_{\tau \geq 1} \sum_{\alpha=1}^{\tau} [I(t, p)]_{\tau, \alpha} t^{\tau} p^{\alpha} + 1 \\
&= \sum_{\tau \geq 1} \sum_{\beta=1}^{\tau} [I(t, p)]_{\tau, \tau+1-\beta} t^{\tau} p^{\tau+1-\beta} + 1 \\
&= \sum_{\tau \geq 1} \sum_{\beta=0}^{\tau} [I(t, p)]_{\tau, \tau+1-\beta} t^{\tau} p^{\tau+1-\beta} + 1 \\
&= 1 + p(\tilde{S}_1(t, p) + 1)
\end{aligned}$$

The end of the proof is straightforward using Proposition 4. \square

Theorem 5.

$$\begin{aligned}
S_1(t, q) &= \tilde{S}_1(t, q^2) \\
&= \frac{1 + (q^2 - 1)t - \sqrt{(q^2 - 1)^2 t^2 - 2(q^2 + 1)t + 1}}{2tq^2}
\end{aligned}$$

5.1.1 Average distance

Theorem 6. *The average edit distance between rooted planar trees with n edges and Id is $n - 1$.*

Proof. 1. The average distance δ can be obtained from the generating function $S_1(t, q)$ in the following way:

- $F(t) = \left. \frac{\partial S_1(t, q)}{\partial q} \right|_{q=1}$
- $\delta = \frac{[F(t)]_n}{C(n)}$ where $C(n)$ is the n -th Catalan number.

This easy computation yields $\delta = n - 1$ but a direct combinatorial interpretation proves this result in a more comprehensive way.

2. This is a direct consequence of Propositions 3 and 4. Another proof can be found in [8, 9]. In [9] the result is more general. Thus we provide here a simpler proof for this special case. \square

5.2 Generating function of the edit distance between one-stack sortable permutations and $n(n - 1) \dots 1$

This is the distance between a tree and the trivial one which is made of n edges and is of height n . It is equivalent to finding the largest decreasing subsequence in the one-stack sortable permutation.

We compute the generating function $D(x, y, z)$ of trees with respect to the number of edges x , the height of the tree y and the number of leaves z at maximal depth.

Proposition 5.

$$D(x, y, z) = yD(x, y, \frac{1}{1 - xz}) - yD(x, y, 1) + \frac{xyz}{1 - xz} \quad (5)$$

Proof.

$$\begin{aligned}
[D(x, y, z)]_{i, j, k} &= \sum_{l=1}^{i-j+1} \binom{l+k-1}{k} [D(x, y, z)]_{i-k, j-1, l} \text{ if } j > 1 \\
[D(x, y, z)]_{i, 1, k} &= \delta_{i, k}
\end{aligned}$$

The coefficient $[D(x, y, z)]_{i,j,k}$ is equal to the number of ways to add k leaves at depth j to any tree with $i - k$ edges, depth $j - 1$ and l leaves at depth $j - 1$. $\binom{l+k-1}{k}$ is the number of ways to add k leaves to l leaves at depth j .

$$\begin{aligned}
D(x, y, z) &= \sum_{i \geq 1} \sum_{j \geq 1} \sum_{k \geq 1} d_{i,j,k} x^i y^j z^k \\
&= \sum_{i \geq 1} \sum_{j \geq 2} \sum_{k \geq 1} \sum_{l \geq 1} \binom{k+l-1}{k} d_{i-k,j-1,l} x^i y^j z^k + y \sum_{i \geq 1} (xz)^i \\
&= \sum_{i \geq 1} \sum_{j \geq 2} \sum_{k \geq 1} \sum_{l \geq 1} (-1)^k \binom{-l}{k} d_{i-k,j-1,l} x^i y^j z^k + y \sum_{i \geq 1} (xz)^i \\
&= \sum_{i \geq 1} \sum_{j \geq 2} \sum_{k \geq 1} \sum_{l \geq 1} (-1)^k \binom{-l}{k} d_{i,j-1,l} x^i y^j (xz)^k + y \sum_{i \geq 1} (xz)^i
\end{aligned}$$

Using

$$(x+a)^{-n} = \sum_{k=0}^{\infty} \binom{-n}{k} x^k a^{-n-k}$$

$$\begin{aligned}
D(x, y, z) &= \sum_{i \geq 1} \sum_{j \geq 2} \sum_{l \geq 1} ((1-zx)^{-l} - 1) d_{i,j-1,l} x^i y^j + y \sum_{i \geq 1} (xz)^i \\
&= y D(x, y, \frac{1}{1-xz}) - y D(x, y, 1) + \frac{xyz}{1-xz}
\end{aligned}$$

□

Let $S_2(x, y)$ be the generating function with respect to the length n of the one-stack sortable permutation and the edit distance between this one-stack sortable permutation and $n(n-1)(n-2) \dots 1$. Then, $S_2(x, y) = D(xy^2, \frac{1}{y^2}, 1)$.

In [10, 11], they give a solution for $D(x, y, 1)$ in terms of a continued fraction.



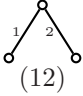

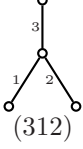
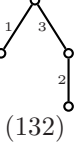
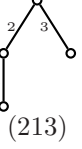
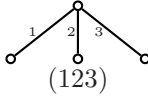
$$D(x, y, 1) = \sum D_k(y) x^k, D_k(y) = \frac{1}{k \left\{ 1 - \frac{y}{1 - \frac{y}{1 - \dots}} \right\}}$$

This yields the solution for S_2 .

$$S_2(x, y) = \sum y^{2k} D_k(\frac{1}{y^2}) x^k$$

The first terms of S_2 are given by:

$$S_2(x, y) = x + x^2 y^2 + x^2 + x^3 y^4 + 3 x^3 y^2 + x^3 + x^4 y^6 + 7 x^4 y^4 + 5 x^4 y^2 + x^4$$

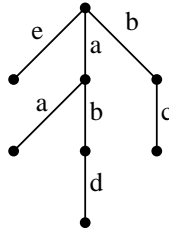
 x (1)		
 x^2 (21)	 (12)	x^2y^2
 x^3 (321)	 (312)  (132)  (213)	 (123) x^3y^4

Average edit distance In [10], they determine analytically the average height of a planar tree with n edges which is $\sqrt{\pi n} - \frac{1}{2}$. Thus, the average edit distance is $2(n - \sqrt{\pi n} + \frac{1}{2}) \equiv 2n$.

6 Conclusion

In section 2.2, we define the edit operations to be insertion and deletion. Indeed we omitted a third one, the relabeling operation. Instead of working with unlabeled trees, we study trees whose vertices are labeled and the relabeling operation consists in changing the label of a vertex.

The general case where the trees are labeled and the different edit operations have different costs can be obtained in a similar way. Define a decorated one-stack sortable permutation as a one-stack sortable permutation where each number is indexed by a letter; $1_e 5_a 2_a 4_b 3_d 7_b 6_c$ represents the following tree:



The operations on decorated one-stack sortable permutations are almost the same as before and the relabeling operation consists in changing one letter. c_i, c_d, c_r are respectively the insert, delete and relabeling unitary costs. There exists only a difference for the insertion of a new free edge. In the unlabeled case, we did not take into account the insertion of a leaf with no sibling. Thus we define a fourth insertion operation as:

- $(\Lambda \xrightarrow{1} i)$ where i is a complete factor of size 1 of the permutation $\sigma = uiv$. $\sigma' = [u]_a[i]_a a[v]_a$ where $a = i$.

Let σ_1 and σ_2 be two decorated one-stack sortable permutations with the same underlying permutation. The label distance $d(\sigma_1, \sigma_2)$ is equal to the string distance between both labeled words.

Let T_1 and T_2 be two decorated one-stack sortable permutations. We denote by a subpermutation σ of T_1 and T_2 a normalized subpermutation without label. Σ_{T_1} is the set of all sub-decorated one-stack sortable permutations of T_1 which underlying permutation is σ .

The relabeling distance between T_1 and T_2 with respect to σ is:

$$d_\sigma(T_1, T_2) = \min\{c_r d(\alpha, \beta), \forall \alpha \in \Sigma_{T_1}, \beta \in \Sigma_{T_2}\}$$

The distance between these two decorated one-stack sortable permutations T_1 and T_2 is given by $\min\{c_i(|T_1| - |\sigma|) + c_d(|T_2| - |\sigma|) + d_\sigma(T_1, T_2), \sigma \text{ normalized subpermutation of } T_1, T_2\}$

References

- [1] K. Zhang and D. Shasha. Simple fast algorithms for the editing distance between trees and related problems. *SIAM J. Comput.*, 18(6):1245–1262, Dec. 1989.
- [2] P.N. Klein. Computing the edit-distance between unrooted ordered trees. In *ESA '98*, pages 91–102, 1998.
- [3] M. Garofalakis and A. Kumar. Correlating XML data streams using tree-edit distance embeddings. In *Proc. PODS'03*, 2003.
- [4] M. Bousquet-Mélou. Sorted and/or sortable permutations. *Disc. Math.*, 225:25–50, 2000.
- [5] J. West. *Permutations and restricted subsequences and Stack-sortable permutations*. PhD thesis, M.I.T., 1990.
- [6] P. Bose, J.F. Buss, and A. Lubiw. Pattern matching for permutations. *Inf. Proc. Letters*, 65:277–283, 1998.
- [7] D.E. Knuth. *The Art of Computer Programming : Fundamental Algorithms*, page 533. Addison-Wesley, 1973.
- [8] E. Deutsch, A.J. Hildebrand, and H.S. Wilf. Longest increasing subsequences in pattern-restricted permutations. *Elect. J. Combin.*, 9(2):R12, 2003.
- [9] A. Reifegerste. On the diagram of 132-avoiding permutations. Technical Report 0208006, Math. CO, 2002.
- [10] N.G. De Bruijn, D.E. Knuth, and S.O. Rice. *Graph theory and Computation*, chapter The average height of planted plane trees. Academic Press, 1972.
- [11] E. Roblet and X.G. Viennot. Théorie combinatoire des t-fractions et approximants de Padé en deux points. *Disc. Math.*, 153:271–288, 1996.